

T cell quantification from DNA sequencing predicts immunotherapy response

Robert Bentham^{1,2*}, Kevin Litchfield^{2,3*}, Thomas B K Watkins^{4*}, Emilia L Lim^{2,4}, Rachel Rosenthal⁴, Carlos Martínez-Ruiz^{1,2}, Crispin T Hiley^{2,4}, Maise Al Bakir⁴, Roberto Salgado^{5,6}, David A Moore^{2,7,8}, Mariam Jamal-Hanjani^{2,8,9}, TRACERx Consortium¹⁰, Charles Swanton^{2,4,8} and Nicholas McGranahan^{1,2+}

1. Cancer Genome Evolution Research Group, Cancer Research UK Lung Cancer Centre of Excellence, University College London Cancer Institute, Paul O’Gorman Building, 72 Huntley Street, London, WC1E 6BT, UK.

2. Cancer Research UK Lung Cancer Centre of Excellence, University College London Cancer Institute, Paul O’Gorman Building, 72 Huntley Street, London, WC1E 6BT, UK.

3. The Tumour Immunogenomics and Immunosurveillance Lab, Cancer Research UK Lung Cancer Centre of Excellence, University College London Cancer Institute, Paul O’Gorman Building, 72 Huntley Street, London, WC1E 6BT, UK.

4. Cancer Evolution and Genome Instability Laboratory, The Francis Crick Institute, 1 Midland Rd, London NW1 1AT, UK

5. Department of Pathology, GZA-ZNA, Antwerp, Belgium

6. Division of Research, Peter MacCallum Cancer Centre, University of Melbourne, Melbourne, Victoria, Australia

7. Department of Cellular Pathology, University College London Hospitals, London, UK

8. Department of Medical Oncology, University College London Hospitals, 235 Euston Rd, Fitzrovia, London, United Kingdom, NW1 2BU, UK

9. Cancer Metastasis Lab, University College London Cancer Institute, Paul O’Gorman Building, 72 Huntley Street, London, WC1E 6BT, UK

10. A list of authors and their affiliations appears at the end of the paper

* These authors contributed equally

+ To whom correspondence should be addressed:

nicholas.mcgranahan.10@ucl.ac.uk

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56

Abstract

The immune microenvironment influences tumour evolution and can be both prognostic and predict response to immunotherapy^{1,2}. However, measuring tumour infiltrating lymphocytes (TILs) is restricted by lack of appropriate data. Whole exome sequencing (WES) of DNA is frequently performed to calculate tumour mutational burden and identify actionable mutations. Here we develop a method for T cell fraction estimation from WES samples, utilising a signal from T cell receptor excision circle (TRECs) loss during VDJ recombination of the T cell receptor alpha (*TCRA*) gene. This score significantly correlates with orthogonal TIL estimates and is agnostic to sample type. Blood *TCRA* T cell fraction is higher in females and correlates with both tumour immune infiltrate and presence of bacterial sequencing reads. Tumour *TCRA* T cell fraction is prognostic in lung adenocarcinoma and using a meta-analysis of immunotherapy-treated tumours, we show that this score predicts immunotherapy response, providing value beyond tumour mutational burden. Applying this score to a multi-sample pan-cancer cohort revealed high diversity in immune infiltration within tumours. Subclonal loss of 12q24.31-32, encompassing *SPPL3*, was associated with reduced *TCRA* T cell fraction. Our method, T cell ExtRECT (T cell Exome TREC Tool), quantifies the T cell infiltrate of WES samples.

Introduction

Checkpoint inhibitors (CPIs) have emerged as revolutionary cancer treatments, acting to release the brakes on the immune system^{3,4}. Clinical response, however, is not universal⁵ and is principally governed by the presence of an immune stimulus, such as neoantigens, and an immune response, mediated by T cells². While neoantigens can be predicted from WES¹, Until now, T cell quantification has required additional biological material, time, and expertise, adding to the cost of immunotherapy.

57 Here we propose a method for the estimation of the T cell fraction present in a WES sample.
58 This method utilises a somatic copy number-based signal from VDJ recombination and the
59 loss of TRECs. We explore the underlying features which predict T cell infiltrate in tumours
60 and blood and evaluate determinants of immune heterogeneity within tumours. Finally, we
61 demonstrate that our estimated T cell fraction can be used as a predictor of clinical response
62 to CPI therapy.

63

64 **Results**

65 ***Inferring T cell fraction from WES data***

66 T cell diversity, which is required for immune system recognition of foreign antigens, is a
67 product of VDJ recombination, where segments within the T cell receptor genes recombine.
68 The alpha chain of the T cell receptor is encoded by the *TCRA* gene (also known as *TRA*) and
69 the result of VDJ recombination is the excision of unselected gene segments from *TCRA* as
70 TRECs, with the T cell undergoing a deletion event within *TCRA*.

71

72 Tools to infer cancer somatic copy number alteration (SCNA)⁶⁻⁹ rely on the read depth ratio
73 (RDR), reflecting the log of the ratio of reads between the tumour sample and its matched
74 control (e.g. buffy coat in a centrifuged blood sample). Deviation in the RDR from zero is
75 assumed to reflect a tumour SCNA. However, within *TCRA* this assumption does not hold; a
76 deviance in the RDR may reflect T cell specific deletion events and SCNA tools may thus
77 erroneously infer tumour SCNA. Indeed, in the TRACERx100 cohort multiple SCNA within
78 *TCRA* were inferred in 165/327 tumour regions (Extended Data Fig. 1a). The RDR deviated
79 the most within segments frequently included within TRECs (Extended Data Fig. 1b-c). This
80 suggests that most detected SCNAs within *TCRA* reflect a signal of relative T cell content
81 rather than cancer SCNAs.

82

83 To exploit this signal to quantify T cell content we developed T cell EXTRECT (T cell Exome
84 TREC Tool). T cell EXTRECT uses a modified RDR within *TCRA* to directly quantify T cell
85 infiltrate in WES samples (Figure 1a), referred to as the *TCRA* T cell fraction. Unlike RNA-seq
86 scores, the *TCRA* T cell fraction represents a direct quantification of the proportion of T cells
87 within a sample. We identified no systematic significant difference in *TCRA* T cell fraction
88 dependent on whether samples were fresh frozen or formalin-fixed paraffin-embedded (FFPE)
89 (Methods, Extended Data Fig. 1d-e). T cell EXTRECT can be applied to any WES sample,
90 thus permitting analysis of T cell fraction in both tumour and blood samples.

91

92

93 ***Validation of TCRA T cell fraction***

94 To evaluate the accuracy of T cell EXTRECT, we used five orthogonal approaches.

95

96 First, to assess the ability to accurately determine the presence or absence of T cells within a
97 sample, we used WES data from cell lines originating from T cell lymphoma (JURKAT, PEER,
98 and HPB-ALL) and 14 colorectal cancer cell lines derived from HCT116 with varying degrees
99 of genomic complexity^{10,11}. All HCT116 cell lines had a calculated fraction of 0. Conversely,
100 the three T cell lymphoma-derived cell lines had scores close to 1 (~0.95-0.96) (Extended
101 Data Fig. 1f).

102

103 Second, we used an alternative DNA based method of inferring immune content¹², based on
104 the number of reads that align to the CDR3 region following VDJ recombination (CDR3 VDJ
105 score, Methods). In the TRACERx100¹³ cohort (Extended Data Fig. 1g) we observed a
106 significant positive correlation between *TCRA* T cell fraction and the CDR3 VDJ score
107 (Extended Data Fig. 1h, $\rho = 0.36$, $P = 1.4e-13$). However, the CDR3 VDJ score was
108 constrained by sequencing depth; the number of reads aligning to the CDR3 region was
109 typically very low (1st quartile = 0, medium = 2, mean = 2.335, 3rd quartile = 3, maximum =
110 14).

111

112 Third, we simulated NGS data with a range of T cell fractions (Extended Data Fig. 2a-d). We
113 observed a highly significant relationship between simulated and calculated T cell fraction (ρ
114 = 0.99986, $P < 2.2e-16$, Extended Data Fig. 2b). Using downsampling and simulations, we
115 found that the *TCRA* T cell fraction estimates remained consistent at coverage above and
116 including 30X ($\rho = 0.84$, $P = 1.4e-14$) (Extended Data Fig. 2e-f). In contrast, the results from
117 the CDR3 method were heavily skewed by sequencing coverage; when selecting the five
118 samples with the highest CDR3 coverage and downsampling to 50X, only one sample with ≥ 3
119 CDR3 reads was detected (Extended Data Fig. 2g).

120

121 Fourth, to further confirm the accuracy of the *TCRA* T cell fraction, we evaluated its association
122 with histopathology-derived TIL scores from H&E slides. Selecting the subset of tumour
123 regions with both RNA-seq data and histopathology-derived TIL scores (147 regions), we
124 evaluated how the *TCRA* T cell fraction, CDR3 VDJ score, and six RNA-seq based immune
125 measures for CD8+ cells (Danaher¹⁴, Davoli¹⁵, xCell¹⁶, TIMER¹⁷, CIBERSORT¹⁸, and EPIC¹⁹)
126 compared to histopathology-derived TIL scores (Figure 1b). The Danaher CD8+ score had
127 the strongest association ($\rho = 0.49$), followed by the *TCRA* T cell fraction ($\rho = 0.41$), Davoli (ρ
128 = 0.4), xCell ($\rho = 0.36$), CIBERSORT ($\rho = 0.23$), TIMER ($\rho = 0.2$), CDR3 VDJ score ($\rho = 0.2$),
129 and EPIC ($\rho = 0.082$).

130

131 Finally, the *TCRA* T cell fraction from WES was directly compared with RNA-seq methods and
132 was found to have a significant positive relationship with multiple immune scores^{1,14-19} with
133 the strongest associations being with T cell related scores (Figure 1c).

134

135 ***Determinants of T cell content in blood***

136 We next explored the key determinants of T cell immune infiltrate in matched control blood
137 WES samples.

138

139 Within the TRACERx100¹³ cohort, blood *TCRA* T cell fraction was significantly higher in
140 females than males (Figure 2a, $P = 0.0057$, $ES = 0.28$) and we observed a trend for higher
141 blood T cell fraction in LUSC compared to LUAD patients (Extended Data Fig. 3a, $P = 0.066$,
142 $ES = 0.19$). We also observed a significant positive relationship between blood *TCRA* T cell
143 fraction and matched tumour *TCRA* T cell fraction (Figure 2a, $\rho = 0.42$, $P = 1.7e-05$). These
144 data suggest that tumour immune infiltrate may influence T cell levels in circulating blood or
145 vice versa. We observed broadly consistent results in LUAD and LUSC TCGA^{20,21} patients
146 (Extended Data Fig. 3b-c).

147

148 To further examine the determinants of blood T cell fraction, we explored WES samples
149 derived from both blood and physiologically normal oesophagus epithelia (PNE) tissue²².
150 While blood samples exhibited a wide range of *TCRA* T cell fraction levels, the majority of
151 PNE tissue had no detectable T cell infiltration (Extended Data Fig. 3d-e). Dividing the PNE
152 samples by presence of T cell infiltration revealed a significant association with blood *TCRA*
153 T cell fraction (Figure 2b, $P = 0.021$, $ES = 0.29$). Therefore, similarly to tumour samples, high
154 levels of T cell infiltration in normal tissue may influence the *TCRA* T cell fraction observed in
155 blood. In a linear model predicting T cell fraction in blood, only the infiltration level in normal
156 tissue was significant (Extended Data Fig. 3f); no genomic factors, such as mutation burden
157 or driver mutation status were predictive of T cell infiltration in PNE tissue (Extended Data Fig.
158 3g).

159

160 Viral or bacterial infections could also influence T cell levels in blood. To explore this we
161 obtained normalised estimates for the abundance of microbial reads from blood and tumour
162 samples from the LUAD and LUSC TCGA cohorts²³. Blood samples with elevated microbial
163 reads ($>$ median, 6.81) had significantly higher blood *TCRA* T cell levels (Figure 2c, $P =$

164 0.00092, ES = 0.31, Wilcoxon test). No corresponding association was identified in tumour
165 samples (Extended Data Fig. 3h, P = NS). No specific virus or bacteria were associated with
166 blood *TCRA* T cell fraction. In tumour samples significant associations for bacteria of the
167 genus *Williamsia* in LUAD ($\rho = -0.17$, P = 0.00011, FDR P = 0.013) and *Paeniclostridium* in
168 LUSC ($\rho = -0.2$, P = 0.00013, FDR P = 0.015) were observed (Extended Data Fig. 3i-k). Both
169 had higher normalised log-cpm values when *TCRA* T cell fraction was lower, suggesting they
170 may be opportunistic species exploiting an immune-cold tumour microenvironment.

171

172 ***Determinants of tumour T cell content***

173 Next, we explored factors influencing T cell infiltrate in tumour tissue. We utilised a recently
174 published pan-cancer cohort of multi-sample data²⁴ to investigate both the extent and possible
175 genomic basis for immune infiltrate heterogeneity. In total, we evaluated T cell infiltrate in 731
176 tumour samples from 178 tumours, from 12 cancer types (Extended Data Fig. 4a-b).

177

178 We classified each multi-sample tumour as uniformly hot (all samples ≥ 0.11 , the mean *TCRA* T
179 cell fraction), uniformly cold (all samples < 0.11) or heterogeneous. There was a significant
180 difference in the proportion of these categories by cancer type (Figure 2d, chi-squared test: P
181 = 1.62e-07) with ER+ breast cancer (BRCA ER+) tumours being the most heterogeneous (83%)
182 and LUSC tumours being the least (22%). Clear differences in the prevalence and heterogeneity
183 of immune infiltrate was observed across cancer types; for instance, while bladder cancer
184 (BLCA) and LUAD had similar numbers of heterogeneous tumours (36% vs 37%), ~64% of BLCA
185 tumours were uniformly immune-hot and 0% were uniformly immune-cold, whereas in LUAD
186 37% tumours were uniformly immune-cold and 25% uniformly immune-hot. This suggests that

187 for certain cancer types there is a highly localised immune infiltrate, which can be subject to
188 considerable sampling bias.

189

190 Next, we examined the relationship between SCNAs and immune diversity. We restricted the
191 analysis to tumours with at least three samples and a heterogeneous mixture of T cell infiltrate.
192 Pairwise SCNA heterogeneity between any two samples was calculated as the sum of the
193 proportion of the genome with unique SCNAs in either region. Pairs of tumour samples with a
194 large disparity in *TCRA* T cell fraction (\geq the mean of all pairwise distances, 0.065) were
195 associated with a larger differences in SCNA heterogeneity compared to matched region pairs
196 with low *TCRA* T cell fraction heterogeneity (Figure 2e, All events: $P = 0.0025$, $ES = 0.347$;
197 gain events: $P = 0.0056$, $ES = 0.318$; loss or LOH events: $P = 0.028$, $ES = 0.253$, $n = 76$).

198

199 To explore whether any specific subclonal SCNA were associated with immune depletion or
200 activation, we identified cytobands that were subclonally lost or gained > 30 tumours in the
201 pan-cancer multi-sample cohort (Extended Data Fig. 4c) and investigated whether specific
202 SCNAs were associated with changes in *TCRA* T cell fraction. Subclonal loss of 12q24.31-32
203 was found to be significantly associated with decreased *TCRA* T cell fraction (Figure 2f: $P =$
204 $5.9e-06$, $ES = 0.75$).

205

206 RNA-seq analysis of the TRACERx100 cohort identified *SPPL3* as exhibiting the most
207 significant differential expression between samples with and without subclonal 12q24.31-32
208 loss (Extended Data Fig. 4d). The absence of *SPPL3* has been found to augment B3GNT5
209 enzyme activity which upregulates cell surface glycosphingolipids that in turn impede class I
210 HLA function and diminish CD8+ T cell activation²⁵. Thus, these data suggest that subclonal
211 loss of 12q24.31, encompassing *SPPL3*, may be selected in tumour evolution across cancer
212 types (occurring in 18.7% of tumours within the cohort) as a mechanism of immune evasion.

213

214 ***T cell fraction is prognostic in LUAD***

215 To explore the clinical utility of T cell ExTRECT, we considered whether the *TCRA* T cell
216 fraction was prognostic in the TRACERx100 non-small cell lung cancer (NSCLC) cohort¹³. We
217 categorised tumour regions as either 'hot' or 'cold' depending on whether *TCRA* T cell fraction
218 was \geq the mean in the cohort (0.081). In LUAD, we observed that patients harbouring an
219 elevated number of immune-cold tumour regions were associated with significantly inferior
220 prognosis (Figure 3, LUAD: ≥ 2 immune-cold regions, HR = 3.1, P = 0.0063 log-rank test,
221 LUAD: ≥ 3 immune-cold regions HR = 7.3, P = 0.00024 log-rank test). In contrast, in LUSC
222 patients there was no significant difference in survival. Using the median (0.074) as a threshold
223 for immune hot or cold regions yielded similar results (Extended Data Fig. 5a). These results
224 are consistent with previous analysis based on TIL scores inferred from computational
225 pathology on the TRACERx100 cohort²⁶. An association between high *TCRA* T cell fraction
226 and good outcome was also observed in the TCGA LUAD (Extended Data Fig. 5b overall
227 survival (OS): HR = 0.61, P = 0.0043, progression free survival (PFS): HR = 0.67 P = 0.016),
228 but not LUSC cohort (Extended Data Fig. 5c). A range of possible thresholds yielded similar
229 results (Extended Data Fig. 5d).

230

231 Consistent with the importance of the tumour region with the lowest immune infiltrate²⁶, the
232 minimum, but not the maximum or mean, *TCRA* fraction across tumour regions was prognostic
233 in the TRACERx100 cohort. Other continuous measures such as a *TCRA* T cell fraction
234 divergence between tumour region score (Extended Data Fig. 5d, LUAD: HR = 2.2 P = 0.023
235 log-rank test) and a model combining both the minimum and maximum scores (Extended Data
236 Fig. 5e, LUAD and LUSC: minimum HR = 0.5, P = 0.005, maximum HR = 1.5 P = 0.061;
237 LUAD: minimum HR = 0.36, P = 0.016, maximum HR = 2.52, P = 0.029) reached significance,

238 suggesting that there is added predictive potential when considering the heterogeneity of the
239 *TCRA* T cell fraction.

240

241 ***T cell fraction and response to CPIs***

242 To further explore the clinical utility of T cell EXTRECT, we evaluated its ability to predict
243 clinical response to CPIs. The CPI1000+ cohort² consists of 1070 CPI-treated tumours
244 receiving either anti-CTLA-4, anti-PD-L1 or anti-PD-1 therapy across eight main cancer types
245 (Extended Data Fig. 6a-b). A responder was defined as a patient with complete response (CR)
246 or partial response (PR), while a non-responder was defined as stable disease (SD) or
247 progressive disease (PD), on imaging by RECIST criteria²⁷.

248

249 Consistent with the importance of T cells in influencing response to CPIs, we observed a
250 significantly higher (Figure 4a, $P = 2.3e-07$, $ES = 0.17$) tumour *TCRA* T cell fraction in
251 responders. Likewise, immune-cold tumours (tumours with *TCRA* T cell fraction < 0.067 , the
252 mean *TCRA* T cell fraction), were significantly enriched for non-responders (Figure 4b,
253 Fisher's exact test, odds ratio (OR) = 2.12, $P = 2.25e-06$).

254

255 Separating the cohort by the medians for both clonal TMB and *TCRA* T cell fraction revealed
256 that the association between *TCRA* T cell fraction and clinical response was independent of
257 clonal TMB (Figure 4b).

258

259 To evaluate the utility of T cell EXTRECT in comparison to RNA-seq based measurements, all
260 studies with ≥ 10 samples from a cancer type with both RNA-seq and *TCRA* T cell fractions
261 were selected for univariate meta-analyses (Figure 4c: 557 patients across 7 studies and 5
262 cancer types). *TCRA* T cell fraction (OR = 1.39, $P = 0.00858$), clonal TMB (OR = 1.59, $P =$
263 $6.021e-05$) and *CD8A* expression (OR = 1.45, $P = 0.0004479$) were all significantly associated
264 with response.

265

266 To assess whether tumour *TCRA* T cell fractions improves prediction of response beyond
267 clonal TMB and to a greater extent than *CD8A* expression we evaluated different linear models
268 (Extended Data Fig. 6c). Only the clonal TMB + *TCRA* model was significant compared to
269 clonal TMB alone (ROC test, $P = 0.0028$, GLM: clonal TMB + *TCRA*, AUC = 0.68, GLM: clonal
270 TMB, AUC = 0.62). When examining the significance of the variables in all models, *TCRA* T
271 cell fraction was more significant than *CD8A* (GLM: clonal TMB + *TCRA*, $P = 4.62e-05$; GLM:
272 clonal TMB + *CD8A*, $P = 0.000431$) and when combined into a multivariable model, *TCRA* T
273 cell fraction remained significant, but *CD8A* expression did not (*TCRA*, $P = 0.00601$, *CD8A*, P
274 $= 0.06246$).

275

276 Finally, we assessed the predictive potential of the *TCRA* T cell fraction in a combined NSCLC
277 CPI cohort (Extended Data Fig. 6d-e) lacking any RNA-seq immune measures. In univariate
278 analyses, (Figure 4d), *TCRA* T cell fraction (OR = 1.44, $P = 0.0071$) and blood *TCRA* T cell
279 fraction (OR = 1.39, $P = 0.015$) were significantly associated with response to CPI. Tumour
280 *TCRA* T cell fraction had OR > 1 in two of three cohorts while blood *TCRA* T cell fraction had
281 OR > 1 in all three cohorts.

282

283 Taken together, these results suggest the *TCRA* T cell fraction can be used as a substitute
284 for RNA-seq measures of CD8+ infiltrate, and, moreover, *TCRA* T cell fraction estimation adds
285 prognostic value to TMB estimates.

286

287 **Discussion**

288 In summary, we present a method, T cell EXTRECT, by which DNA WES can be harnessed
289 to study the immune microenvironment. T cell EXTRECT provides an accurate estimate of
290 immune infiltrate which shows clinical utility. We find tumour *TCRA* T cell fraction is prognostic
291 in LUAD and validate this finding in the TCGA LUAD cohort. Relatedly, we find the *TCRA* T
292 cell fraction is associated with response to CPI in a pan-cancer cohort and improves upon the

293 predictive value of clonal TMB. T cell ExtRECT enables the T cell fraction to be calculated in
294 any WES sample. Leveraging this, we demonstrate that T cell fraction in blood is
295 heterogeneous, associated with microbial infections and was found to be significantly higher
296 in females than males in TRACERx100 NSCLC patient data, consistent with previous
297 findings^{30,31}. Our analysis of blood samples in the lung CPI cohort revealed that blood *TCRA*
298 T cell fraction is predictive of response to immunotherapy.

299

300 The T cell ExtRECT method has limitations. While the tool provides a quantification of the
301 proportion of T cells in a sample, it cannot distinguish neoantigen-reactive from bystander T
302 cells, and is unable to detect clonotypes. Further, T cell ExtRECT loses fidelity below 30X
303 sequencing depth. Nevertheless, this relatively low depth means it should be applicable to
304 most DNA sequencing datasets. T cell ExtRECT has so far only been optimised for WES, but
305 further work will extend the method to whole-genome and to other species including much
306 studied model organisms. T cell ExtRECT has clear applications in the immuno-oncological
307 exploration of tumour samples, however it could also be utilised in a wider clinical setting, such
308 as newborn screening of severe combined immunodeficiency disease³².

309

310 In summary, our approach, T cell ExtRECT, could have important applications in both basic
311 and translational research by providing a cost-effective technique to characterise immune
312 infiltrate alongside somatic changes, without the need for RNA sequencing.

313

314 ***Methods***

315

316 A detailed and full description of the T cell ExtRECT method is given in Supplementary
317 Information.

318

319 ***Statistics***

320 All statistical tests were performed in R 3.6.1. No statistical methods were used to
321 predetermine sample size. Tests involving correlations were done using 'stat_cor' from R
322 package ggpubr (v0.4.0) with the Spearman's method. Tests involving comparisons of
323 distributions were done using 'stat_compare_means' using either 'wilcox.test' using the
324 unpaired option, unless otherwise stated. Effect sizes for the corresponding Wilcoxon tests
325 were measured using the 'wilcox_effsize' function from the rstatix package (v0.6.0). Hazard
326 ratios and P values were calculated with the 'survival' package (v3.2-3) for both the Kaplan-
327 Meier curves and Cox proportional hazard model. For all statistical tests, the number of data
328 points included are plotted or annotated in the corresponding figure. Plotting and analysis in
329 R also made use of the ggplot2 (v3.3.3), dplyr (v1.0.4), tidyr (v1.1.1), gridExtra (v2.3) and
330 gtable (v0.3.0) packages.

331

332 *Fresh frozen vs FFPE samples*

333 To test that the *TCRA* T cell fraction was reliable and consistent for both fresh frozen and
334 FFPE samples the non-GC corrected *TCRA* T cell fractions were calculated for six different
335 studies within the CPI1000+ cohort. Three of these studies utilised WES derived from FFPE
336 tissues (n = 460) while the other three utilised WES samples derived from fresh frozen tissue
337 (n = 357).

338

339 Fitting a linear model to predict *TCRA* T cell fraction by histology and FFPE status (Extended
340 Data Fig. 1i) revealed that cancer type however was the main driver of this significance with
341 FFPE status not being significant. Additionally, for melanoma and bladder tumours that had
342 FFPE and fresh frozen WES samples there was no significant difference found (Extended
343 Data Fig. 1f). This led us to conclude that whether a WES sample is derived from fresh frozen
344 or FFPE tissue does not significantly affect the values of the *TCRA* T cell fraction calculated
345 by T cell ExtRECT.

346

347 *Calculation of CDR3 VDJ scores*

348 The procedure outlined in Levy et al.¹² was followed to calculate the CDR3 VDJ scores. First
349 reads aligning to *TCRB* (hg19:chr7:142000817-142510993) and unaligned reads were
350 extracted with samtools, this resulting bam was converted to fastq using bedtools and then
351 the tool IMSEQ (v1.1.0)³³ was used on the resulting output to identify VDJ recombinant reads
352 aligning to the CDR3 region, the number of aligned reads was then normalised by the total
353 number of reads in the original bam file (as measured by samtools flagstat) to create the CDR3
354 VDJ scores.

355

356 *Kraken TCGA analysis*

357 Pre-processed microbiome data output from the Kraken³⁴ analysis performed by Poore et al.²³
358 was downloaded from ftp://ftp.microbio.me/pub/cancer_microbiome_analysis/.

359

360 To create the high and low Kraken microbiome groups for both the blood and tumour samples
361 the file *Kraken-TCGA-Voom-SNM-Most-Stringent-Filtering-Data.csv* was downloaded
362 containing normalised log-cpm values, for each sample the rows were summed giving a
363 overall 'microbiome' score. The samples were then divided into high and low groups based on
364 the median of this score.

365

366 To investigate the role of any individual microbial species in influencing *TCRA* T cell fraction
367 a reduced list of the species from the *Kraken-TCGA-Voom-SNM-Most-Stringent-Filtering-*
368 *Data.csv* file were selected, by removing all species with less than 1000 total raw reads in the
369 TCGA LUAD and LUSC cohort as called from the raw data file *Kraken-TCGA-Raw-Data-*
370 *17625-Samples.csv*. This left a total of 59 microbial species that were individually tested for
371 association with *TCRA* T cell fraction using Spearman's correlation for both LUAD and LUSC
372 blood and tumour samples.

373

374 *TRACERx100 patients*

375 The first 100 patients prospectively analysed by the NSCLC TRACERx study
376 (<https://clinicaltrials.gov/ct2/show/NCT01888601>, approved by an independent research
377 ethics committee, 13/LO/1546) were used in this study. This is identical to the 100 patient
378 cohort originally described in Jamal-Hanjani et al¹³.

379

380 Describing this cohort in brief, informed consent was a mandatory requirement for entry into
381 the TRACERx study. This NSCLC cohort consisted of 68 males and 32 female patients with
382 a median age of 68. Finally, the cohort is predominantly made up of early-stage tumours (Ia
383 (26), Ib (36), IIa (13), IIb (11), IIIa (13) and IIIb (1)) and 28 patients also had adjuvant therapy.

384

385 *TRACERx100 WES and RNA-seq samples*

386 Both WES (aligned to hg19) and RNA-seq samples were obtained from the TRACERx study
387 for the first 100 patients, the method for processing these samples is as previously
388 described¹³. Notably for the WES samples, exome capture was performed using a custom
389 version of Agilent Human All Exome V5 kit as per the manufacturer instructions.

390

391 *TCGA LUAD and LUSC cohorts*

392 Aligned BAM files (hg38) from the TCGA LUAD and LUSC cohorts were downloaded from the
393 genomic data commons (dataset ID: phs000178.v10.p8). Sample purity and ploidy calls were
394 generated from ASCAT (v2.4.2) from a previous analysis of the TCGA data³⁵, in short
395 Affymetrix SNP6 profiles from paired tumour-normal samples (dataset ID: phs000178.v10.p8)
396 were processed by PennCNV libraries³⁶ to obtain BAFs and log ratios which were GC
397 corrected before being processed with ASCAT⁶.

398

399 *Cancer cell line data*

400 The non-T cell derived colorectal cancer cell lines HCT116 were sequenced with Illumina
401 HiSeq 2500 and aligned with bwa mem using hg19 as described in López et al.¹⁰. The T cell
402 derived cell lines were from the dataset were described in Ghandi et al.¹¹ and downloaded

403 from the Sequence Read Archive (SRA) under accession number PRJNA523380. Cell lines
404 derived from T cells were chosen ensuring that any cell line derived from precursor T cell
405 acute lymphoblastic leukemia were excluded as these have not undergone VDJ
406 recombination. This process led to WES data from three cell lines being chosen: JURKAT,
407 HPB-ALL, and PEER.

408

409 Due to the difficulty of running ASCAT without matching germline samples, the naïve *TCRA* T
410 cell fraction was used for all cell line work.

411

412 *Multi-sample tumour cohort of patients*

413 The multi-sample pan-cancer cohort (Extended data Fig. 4b) was created by combining the
414 TRACERx cohort with a subset of the cohort presented recently by Watkins et al.²⁴. Tumours
415 were included if they had at least two regions sequenced in the primary tumour for which it
416 was possible to calculate the *TCRA* T cell fraction using T cell EXTRECT. The final cohort
417 therefore consisted of a multi-region primary tumour data set with the addition of any
418 metastasis samples that were also sequenced for these patients.

419 Besides TRACERx100 the following datasets were combined into the final multi-sample pan-
420 cancer cohort:

421

- 422 1. Brastianos et al.³⁷ - a cohort focused on studying brain metastasis originating from
423 different histologies, only tumours with multi-region primary samples from this cohort
424 were included.
- 425 2. Gerlinger et al.^{38,39} - A multi-sample primary cohort of renal clear cell carcinoma (KIRC)
426 patients.
- 427 3. Harbst et al.⁴⁰ - A multi-region primary cohort of skin cutaneous melanoma (SKCM)
428 patients.
- 429 4. Lamy et al.⁴¹ - A multi-region primary cohort of bladder cancer patients (BLCA)

- 430 5. Savas et al.⁴² - A multi-sample cohort of ER+ and triple-negative breast cancer patients
431 (BRCA ER+ and TNBC)
- 432 6. Suzuki et al.⁴³ - A multi-region primary cohort of glioma.
- 433 7. Turajlic et al.⁴⁴ - A multi-region primary cohort of clear cell renal cell carcinoma (KIRC).
- 434 8. Messaoudene et al.⁴⁵ - A multi-region primary cohort of HER2+ and ER+ breast cancer
435 patients.

436

437

438 *Selection of subregions for multi-region sequencing in different data sets*

439 In all of the multi-region cohorts regions were selected though by different methods (see
440 associated publications) with two main criteria in mind, first that tumour content be maximised
441 at the expense of stromal in order to assure good quality mutation and copy number analysis
442 for the main goal of the genomic analysis and second that each region represent a physically
443 separate and distinct part of the tumour. In cases where these were not at separate sites
444 different measures were used. In the TRACERx100 cohort for example regions sequenced
445 were a minimum of 3mm apart.

446

447 *Identification of gain, loss, and LOH events in a pan-cancer multi-sample cohort*

448 Analysis of whole-exome sequencing was performed as described previously¹³. Copy-number
449 segmentation, tumour purity and ploidy for each sample were estimated using ASCAT⁶ as
450 described previously¹³. These data were used as input to a multi-sample SCNA estimation
451 approach to produce genome-wide estimates of the presence of loss of heterozygosity as well
452 as loss, neutral, gain and amplification copy-number states relative to sample ploidy. The log
453 ratio values present in each copy-number segment with ≥ 5 log ratio values in all samples of a
454 tumour were examined relative to three sample-ploidy-adjusted log ratio thresholds using one-

455 tailed t-tests with a $P < 0.01$ threshold. These log ratio thresholds were equivalent to
456 $<\log_2[1.5/2]$ for losses, $>\log_2[2.5/2]$ for gains in a diploid tumour. Any segment not classified
457 as a loss or gain were classed as neutral. For each segment, these relative to ploidy definitions
458 were combined with loss of heterozygosity detection across all samples from a single tumour.

459

460 *Pairwise subclonal SCNA scores*

461 To calculate pairwise subclonal SCNA measures, the classifications outlined in the previous
462 methods section were used to create three groups of pairwise subclonal SCNA scores. First,
463 we considered any segment affected by any of gain or loss relative to ploidy or LOH as
464 aberrant and compared each pair of regions from a single patient's disease, classifying
465 aberrant areas as clonal if aberrant in both samples or subclonal if aberrant in only one
466 sample. This same process was repeated for gains relative to ploidy alone and then losses
467 relative to ploidy and LOH considered together.

468

469 *Cytoband-level SCNA analysis*

470 To enable comparisons across tumours, segments were mapped to hg19 cytobands. If
471 multiple segments mapped to a cytoband, the SCNA status (gain or loss relative to ploidy) of
472 the segment with the largest overlap with the cytoband was chosen.

473

474 For the SCNA gain and loss analysis, cytoband level events were selected if they occurred
475 subclonally across the entire cohort greater than 30 times. Bands passing this threshold within
476 the same region (e.g. all cytobands on 1p36) were then grouped together. A Wilcoxon paired
477 test was used to assess whether the tumour regions within a single patient with the subclonal
478 SCNA events had a significant difference in *TCRA* T cell fraction to those regions without the
479 event.

480

481 *Selection of multi-sample tumours with heterogeneous immune infiltration*

482 To be included a tumour had to have at least 3 regions sequenced and meet the following two
483 requirements, 1) have a pair of regions with a large change in immune infiltration as defined as
484 having ≥ 0.065 difference in *TCRA* T cell fraction, and 2) have a pair of regions with a small
485 or no change in immune infiltration as defined as having < 0.065 difference in *TCRA* T cell
486 fraction. An example of a tumour matching this requirement would be one with three regions
487 R1, R2 and R3 with *TCRA* T cell fractions of 0.01, 0.01 and 0.2 respectively. The R1-R2 pair
488 has a difference in *TCRA* T cell fraction of 0 while the R1-R3 and R2-R3 pairs would both
489 have a large difference of 0.19. Within the multi-sample tumour cohort 76 patients matched
490 these criteria.

491

492 *RNA-seq differential gene expression analysis for patients with subclonal 12q24.31-32 loss*

493 Differential gene expression analysis was performed on the TRACERx100 RNA-seq patients
494 with subclonal 12q24.31-32 loss. Using R 4.0.0, first the edgeR R package (version 3.32.1)
495 was used for sample-specific TMM normalisation, any genes with low expression were then
496 filtered out using the standard edgeR filtering method before using the Limma-Voom method
497 from the limma R package (version 3.46.0) to calculate the Voom fit and obtain p-values for
498 the gene expression differences. The comparison controlled for patient and histology as
499 blocking factors and p-values were FDR corrected for multiple testing. Results were then
500 visualised with the R EnhancedVolcano package (version 1.8.0).

501

502 *CPI1000+ meta-analysis of cohorts*

503

504 The CPI1000+ cohort is fully described in Litchfield et al.² and contains the following datasets:

- 505 1. Snyder et al.⁴⁶, an advanced melanoma anti-CTLA-4 treated cohort.
- 506 2. Van Allen et al.⁴⁷, an advanced melanoma anti-CTLA-4 treated cohort.
- 507 3. Hugo et al.⁴⁸, an advanced melanoma anti-PD-1 treated cohort.

- 508 4. Riaz et al.⁴⁹, an advanced melanoma anti-PD-1 treated cohort.
- 509 5. Cristescu et al.⁵⁰, an advanced melanoma anti-PD-1 treated cohort.
- 510 6. Cristescu et al.⁵⁰, an advanced head and neck cancer anti-PD-1 treated cohort.
- 511 7. Cristescu et al.⁵⁰ “all other tumour types” cohort (from KEYNOTE-028 and KEYNOTE-
- 512 012 studies), treated with anti-PD-1.
- 513 8. Snyder et al.⁵¹, a metastatic urothelial cancer anti-PD-L1 treated cohort.
- 514 9. Mariathasan et al.⁵², a metastatic urothelial cancer anti-PD-L1 treated cohort.
- 515 10. McDermott et al.⁵³, a metastatic renal cell carcinoma anti-PD-L1 treated cohort.
- 516 11. Rizvi et al.²⁹, a non-small cell lung cancer anti-PD-1 treated cohort.
- 517 12. Hellman et al., a cohort of non-small cell lung cancer samples treated with anti-PD-1
- 518 used by Litchfield et al.².
- 519 13. Le et al.⁵⁴, a colorectal cancer cohort treated with anti-PD-1 therapy.

520

521 Of these studies Snyder et al.⁵¹ was excluded from the analysis due to extremely poor

522 coverage within the *TCRA* gene. Additionally, 55 patients were either on treatment at the time

523 of the biopsy or had prior treatment with CPIs and were removed from the analysis. All

524 samples were aligned to hg19 using bwa mem (v0.7.15) with purity and SCNA data calculated

525 using ASCAT as described in Litchfield et al.².

526 Notably, 953/1070 (89%) samples had WES data, 888/1070 (83%) had sufficient purity and

527 coverage to enable copy number calculation enabling the *TCRA* T cell fractions to be

528 calculated. 643/1070 (60%) of these samples had matched RNA-seq data allowing orthogonal

529 assessment of T cell estimates.

530

531 For an extension to this dataset, Shim et al.²⁸ a NSCLC anti-PD-1 treated cohort was added

532 for a specific NSCLC analysis. In this entire cohort mutations were called as either clonal or

533 subclonal using PyClone as described by Litchfield et al.².

534

535 *Orthogonal immune measures*

536 *RNA-seq signatures*

537 We used the method of Danaher et al.¹² as our primary method of estimating T cell content
538 from RNA-seq measures as it has been previously demonstrated that this is most strongly
539 correlated to TIL scores calculated in TRACERx¹. Other RNA-seq signatures tested against
540 the *TCRA* T cell fractions were the Davoli method¹⁵, xCell¹⁶, TIMER¹⁷ and EPIC¹⁹ and
541 CIBERSORT¹⁸.

542

543 *Histopathology-derived TIL scores*

544 TILs were estimated, as previously described in Rosenthal et al.¹, from histopathology slides
545 using internationally established guidelines, developed by the International Immuno-Oncology
546 Biomarker Working Group⁵⁵. In brief, the relative proportion of stromal area to tumour area
547 was determined from the pathology slide of a given tumour region. TILs were reported for the
548 stromal compartment (= percent stromal TILs). The denominator used to determine the
549 percent stromal TILs was the area of stromal tissue (that is, the area occupied by mononuclear
550 inflammatory cells over total intratumoral stromal area) rather than the number of stromal cells
551 (that is, the fraction of total stromal nuclei that represent mononuclear inflammatory cell
552 nuclei). This method has been demonstrated to be reproducible among trained pathologists⁵⁶.
553 An inter-person concordance was performed, and this demonstrated high reproducibility. The
554 International Immuno-Oncology Biomarker Working Group has developed a freely available
555 training tool to train pathologists for optimal TIL assessment on haematoxylin–eosin slides
556 (www.tilsincancer.org).

557

558 *Univariate and multivariable model for CPI response*

559 For the univariate model an adapted procedure from Litchfield et al.² was followed with the
560 main difference being that only samples with complete data (RNA-seq for *CD8A*, clonal TMB
561 and *TCRA* T cell fraction) were included. The univariate model meta-analysis was conducted
562 using R package 'meta' (version 4.13-0). The multivariable model was created with general

563 linear models using the function `glm` from the 'stats' R package using default values. The R
564 package 'ROCR' (version 1.0-11) was used for the ROC curve analysis.

565

566 *Code*

567 The code used to produce *TCRA* T cell fraction scores is available for academic non-
568 commercial research purposes upon reasonable request.

569

570 All other code used in the analysis and to produce figures is available at:
571 <https://github.com/McGranahanLab/T-cell-ExtRECT-figure-code-2021>

572

573 *Data availability*

574 The RNA-seq data and WES data (in each case from the TRACERx study) generated, used
575 or analysed during this study are not publicly available and restrictions apply to the availability
576 of these data. Such RNA-seq and WES data are available through the Cancer Research UK
577 & University College London Cancer Trials Centre (ctc.tracerx@ucl.ac.uk) for academic non-
578 commercial research purposes upon reasonable request, and subject to review of a project
579 proposal that will be evaluated by a TRACERx data access committee, entering into an
580 appropriate data access agreement and subject to any applicable ethical approvals.

581

582 Details of all other datasets obtained from third parties used in this study can be found in
583 Extended Data Table 1. Clinical trial information (if applicable) is also available within the
584 associated publications described in Extended Data Table 1.

585

586 **Figure Legends**

587 **Figure 1 – Overview and validation of T cell ExtRECT**

588 **a**, Overview of how VDJ recombination signal is identified from read depth within *TCRA* in T
589 cell fraction calculation. **b**, Association with histopathology TIL scores and measures of CD8+

590 T cell content from either RNA-seq (Danaher, Davoli, EPIC, TIMER, CIBERSORT and xCell)
591 or DNA (T cell ExTRECT and CDR3 VDJ score). **c**, Association between *TCRA* T cell fraction
592 with RNA-based scores for immune cell types (Danaher¹⁴, Davoli¹⁵, EPIC¹⁹, TIMER¹⁷,
593 CIBERSORT¹⁸, and xCell¹⁶) ordering determined by strength of association (Spearman's Rho
594 coefficient) with *TCRA* T cell fraction.

595

596 **Figure 2: Determinants of T cell fraction.**

597 **a**, TRACERx100 blood *TCRA* T cell fraction predictors. **b**, Association of *TCRA* T cell fraction
598 in PNE with blood *TCRA* T cell fraction. **c**, Microbial reads from Kraken versus blood *TCRA* T
599 cell fraction (n = 111). **d**, Proportion of tumours uniformly immune-hot, uniformly immune-cold
600 or heterogeneous (Methods). **e**, Multi-sample tumours (n = 76) with heterogeneous immune
601 infiltrate defined as having both a pair of regions with pairwise *TCRA* T cell fraction difference
602 < 0.065 and another with pairwise difference \geq 0.065, versus pairwise SCNA heterogeneity
603 score (Methods). Threshold 0.065 being the mean of all pairwise differences between regions.
604 **f**, *TCRA* T cell fraction difference between regions with or without subclonal loss of 12q24.31-
605 32. All Wilcoxon tests two sided and boxplots represent lower quartile, median, and upper
606 quartile.

607

608 **Figure 3 – Prognostic value of *TCRA* T cell fraction within LUAD but not LUSC**

609 TRACERx100 multi-region LUAD (top) and LUSC (bottom) Kaplan-Meier curves divided by
610 the number of immune-cold regions in the tumour (increasing left to right). Immune-hot and
611 immune-cold regions defined using threshold of the mean of all tumour regions (0.08095).
612 Patients in Kaplan-Meier analyses were restricted to those with total regions greater than the
613 number of immune-cold regions used in defining the threshold.

614

615 **Figure 4 – *TCRA* T cell fraction is predictive of survival and response to immunotherapy**

616 **a**, Violin plot showing the tumour *TCRA* T cell fraction for non-responders versus responders
617 across the CPI1000+ cohort, dotted black line shows mean *TCRA* T cell fraction (0.067) **b**,
618 Tumour *TCRA* T cell fraction versus clonal TMB, dashed lines divide cohort into four quadrants
619 with high/low clonal TMB and immune-hot/immune-cold tumours separated by the median
620 values. Inset pie charts indicate the percentage of patients demonstrating CPI response. **c**,
621 Univariate meta-analysis of predictors of CPI response across multiple cohorts with ≥ 10
622 patients of a cancer type and both DNA and RNA-seq data. Left panel: forest plot of OR values
623 from different clinical factors with associated p-values in terms of predictive value of response.
624 Right panel: heatmap of OR values across individual studies from the CPI1000+ dataset,
625 focusing on cohorts with both RNA-seq and *TCRA* T cell fraction. **d**, Univariate meta-analysis
626 across three CPI lung datasets with DNA but no RNA- seq data.

627

628 **Extended Data Fig. 1: Overview and validation of T cell ExTRECT**

629 **a**, Outline of quantification of the *TCRA* T cell fraction utilising VDJ recombination and TRECs.
630 *top*: Schematic demonstrating how RDR signals are used to detect SCNA gain or loss events
631 in a standard tumour and matched control sample analysis. In this analysis cells consist of
632 three distinct cell types: tumour cells, T cells and all other stromal cells. *bottom*: Schematic of
633 how this same process works when focussing on the *TCRA* gene in relation to VDJ
634 recombination and TRECs, the lower right panel indicates an increased number of breakpoints
635 detected in the TRACERx100 dataset within the *TCRA* gene relative to surrounding areas of
636 14q, suggesting that the TREC signal is captured. **b**, **c**, Plots showing examples of RDR in
637 two TRACERx100 regions demonstrating either increased levels of T cell content in blood
638 compared to matched tumour (**b**) or increased levels of T cell content in tumour compared to
639 matched blood (**c**). VDV segments refer to variable segments in both the *TCR α* and *TCR δ*
640 locus. **d**, *TCRA* T cell fraction (non-GC corrected) value for FFPE and fresh frozen samples

641 for bladder and melanoma tumours within the CPI1000+ cohort (bladder: n = 228, melanoma:
642 n= 297, two sided wilcoxon test used, boxplot shows lower quartile, median and upper quartile
643 values). **e**, Summary of linear model for prediction of non-GC corrected *TCRA* T cell fraction
644 from histology and FFPE sample status within the CPI cohort. **f**, Pie charts of calculated *TCRA*
645 T cell fraction from WES of either T cell-derived cell lines or non-T cell derived cell lines, all
646 HCT116 cell lines had calculated fractions < 1 e-15. **g**, Overview of samples in the
647 TRACERx100 cohort. **e**, Association of the CDR3 VDJ read score based on the iDNA method
648 to *TCRA* T cell fraction in TRACERx100, error bands represent the 95% confidence interval
649 of the fitted linear model.

650

651 **Extended Data Fig. 2: Accuracy of *TCRA* T cell fraction by copy number and depth**

652 **a**, Simulated log RDR from a sample consisting of 24% T cells, 75% tumour, and 1% non-T
653 cell stroma (*TCRA* copy number = 1). **b**, Calculated *TCRA* T cell fraction versus actual T cell
654 fraction value for simulated data **c**, Difference between calculated naïve T cell fraction and
655 actual fraction for range of tumour purities and local tumour copy number states at the *TCRA*
656 locus. **d**, Difference between *TCRA* T cell fraction and actual fraction for a range of tumour
657 copy number and purities. **e**,. Downsampling of 5 TRACERx100 regions to different depths. **f**,
658 Downsampling of simulated data to different depth levels. **g**, Downsampling of the 5
659 TRACERx100 regions that with the highest CDR3 read counts to different depths and the
660 resulting CDR3 read counts.

661

662 **Extended Data Fig. 3: Extended analysis on determinants of *TCRA* T cell fraction**

663 **a**, Association of blood *TCRA* T cell fraction to histology in TRACERx100 (n = 93 LUAD and
664 LUSC patients, two sided wilcoxon test used for P value). **b**, Predictors of blood *TCRA* T cell
665 fraction in TCGA LUAD and LUSC cohort (left panel: n = 1017, middle panel: n = 976, right
666 panel: n = 714). **c**, Overview of samples in the TCGA LUAD and LUSC cohort. **d**, Summary
667 of mean *TCRA* T cell fraction in PNE cohort. **e**, Overview plot of PNE cohort containing multi-
668 region microdissected tissue paired with normal blood samples. **f**, Summary of linear model

669 for predicting blood *TCRA* T cell fraction, PNE infiltration defined as *TCRA* T cell fraction >
670 0.001, ESCC = Oesophageal squamous cell carcinoma, HGD = high grade dysplasia. **g**,
671 Linear model for *TCRA* T cell fraction in PNE samples from genomic factors. **h**, Association
672 of microbial reads from Kraken with *TCRA* T cell fraction in tumour samples (n = 880). **i**, -
673 Log10 p-values for 59 microbial species tested for association with *TCRA* T cell fraction in
674 blood and tumour sample in LUAD and LUSC. Red line represents the significance threshold
675 at P = 0.000423. **j**, The significant hit *Williamsia* in LUAD tumours, red dots represent samples
676 where reads were detected while blue represent samples with no reads detected (n = 501).
677 **k**, The significant hit *Paeniclostridium* in LUSC tumours (n = 379). All wilcoxon tests two sided
678 and boxplots represent lower quartile, median and upper quartile.

679

680 **Extended Data Fig. 4: Subclonal SCNAs and T cell infiltration**

681 **a**, Overview of immune heterogeneity across multi-sample pan-cancer cohort with tumour
682 regions ranked by *TCRA* T cell fraction, *upper panel*: histogram of entire cohort, *lower panel*:
683 tumour regions grouped by patients with solid horizontal lines joining regions from the same
684 patient, each line includes 2 or more tumour region and dashed red line is at the mean *TCRA*
685 T cell fraction in the cohort (0.11). **b**, Overview of patients in the multi-sample pan-cancer
686 cohort. **c**, Lower panel: number of tumours in pan-cancer multi-sample cohort with subclonal
687 gains (dark red) or losses (dark blue) across the genome, horizontal lines signify the regions
688 which have more than 30 tumours (Methods) with subclonal gains or losses. *Upper panel*: -
689 log10(p-value) of the 160 cytoband regions tested for association between *TCRA* T cell fraction
690 and subclonal gains (dark red points) or losses (dark blue points). Red horizontal line marks
691 significance threshold, only one region is significant, a loss event on chromosome 12q24.31-
692 32. **d**, Volcano plot for the RNA-seq analysis in the TRACERx100 cohort between regions with
693 12q24.31-32 loss and regions without, genes within the locus are labeled, dotted lines at fold
694 change of 0.25 and adjusted P = 0.05.

695

696 **Extended Data Fig 5 : Association of *TCRA* T cell fraction with prognosis**

697 **a**, Kaplan-Meier curves for the multi-region TRACERx100 cohort for LUAD (top) and LUSC
698 (bottom) divided by the number of cold regions in the tumour. Hot and cold regions were
699 defined by using the median of all the tumour regions (0.0736) as a threshold. In each Kaplan-
700 Meier curve the included patients were restricted to those with total regions greater than the
701 number of cold regions used in defining the threshold. **b**, Kaplan-Meier curves for overall and
702 progression free survival in the TCGA LUAD cohort, dividing the cohort into immune hot and
703 cold groups using the mean of the TCGA LUAD cohort (0.109) as a threshold. **c**, Kaplan-
704 Meier curves for the TCGA LUSC, and TCGA LUAD & LUSC cohorts for overall and
705 progression free survival using the mean of the TCGA LUAD cohort (0.109) as a threshold for
706 distinguishing hot and cold tumours. **d**, Log₂(Hazard ratios) from Kaplan-Meier plots for the
707 TCGA separating the tumour samples into hot and cold based on different thresholds from 0
708 to 0.16 in steps of 0.0025 for overall and progression free survival. **e**, Hazard ratios of separate
709 Cox regression models relating disease free survival to different multi-region measures
710 related to the *TCRA* T cell fraction in the entire TRACERx100 cohort as well as the LUAD and
711 LUSC patients separately. *TCRA* divergence score is defined as the maximum divided by the
712 upper 95% confidence interval of the minimum. **f**, Hazard ratios of separate Cox regression
713 models for *TCRA* T cell fraction for the TCGA LUAD and LUSC cohort for both overall survival
714 (OS) and progression free survival (PFS).

715

716 **Extended Data Fig 6: Overview of CPI1000+ cohort**

717 **a**, Cohort overview of the CPI1000+ dataset. **b**, Overview of samples in the CPI1000+ cohort
718 excluding Snyder et al., 2017 and those with prior CPI treatment. **c**, ROC plot of GLM models
719 for predicting CPI response (blue: clonal TMB, red: clonal TMB + *TCRA* T cell fraction, green:
720 clonal TMB + *CD8A* expression). **d**, Cohort overview of the CPI lung dataset, red lines in
721 upper panel reflect the median *TCRA* T cell fraction in patients with (0.10) or without (0.0070)
722 a response to CPI, note that Tumour *TCRA* T cell fraction particularly in non-responders is
723 often zero. **e**, Overview of patients in the CPI Lung cohort.

724

725 **Extended Data Table 1: Original source publications**

726 Original source publications (excluding TRACERx studies) containing the sequencing data
727 used in either the multi-sample pan-cancer cohort, PNE cohort or the CPI1000+ cohort.
728 Studies including lung cancer patients used in the lung CPI cohort are noted.

729

730

731

732

733 **Author contributions**

734 R.B. helped conceive the study, designed and conducted the bioinformatic analysis, and wrote
735 the manuscript. K.L. curated the CPI1000+ cohort used in the study and provided considerable
736 bioinformatic support on its analysis. T.B.K.W. provided considerable bioinformatic support on
737 the analysis of the multi-sample pan-cancer cohort and helped conceive the study and write
738 the manuscript. T.B.K.W. and E.L.L jointly curated the multi-sample pan-cancer cohort used
739 in the study. R.R. and C.M.-R. provided considerable bioinformatic support in the
740 transcriptomic analysis performed in the study, providing RNA-seq immune score metrics and
741 assisting with the RNA-seq gene expression analysis respectively. R.S., M.A.B., D.A.M., and
742 C.T.H. jointly analysed histopathology-derived TIL estimates. C.S. helped provide study
743 supervision and helped direct the avenues of bioinformatics analysis and also gave feedback
744 on the manuscript. N.M conceived and supervised the study and helped write the manuscript.

745

746 **Competing interests**

747 D.A.M. reports speaker fees from AstraZeneca. M.A.B. has consulted for Achilles
748 Therapeutics. R.R. has consulted for and has stock options in Achilles Therapeutics. K.L.
749 reports speaker fees from Roche Tissue Diagnostics. C.T.H. has received speaker fees from
750 AstraZeneca. M.J.-H. is a member of the Scientific Advisory Board and Steering Committee
751 for Achilles Therapeutics. N.M. has stock options in and has consulted for Achilles
752 Therapeutics and holds a European patent in determining HLA LOH (PCT/GB2018/052004).

753 C.S. acknowledges grant support from Pfizer, AstraZeneca, Bristol Myers Squibb, Roche-
754 Ventana, Boehringer-Ingelheim, Archer Dx Inc. (collaboration in minimal residual disease
755 sequencing technologies) and Ono Pharmaceutical; is an AstraZeneca Advisory Board
756 Member and Chief Investigator for the MeRmaiD1 clinical trial; has consulted for Amgen,
757 Pfizer, Novartis, GlaxoSmithKline, MSD, Bristol Myers Squibb, AstraZeneca, Illumina,
758 Genentech, Roche-Ventana, GRAIL, Medicxi, Bicycle Therapeutics, Metabomed and the
759 Sarah Cannon Research Institute; has stock options in Apogen Biotechnologies, Epic
760 Bioscience and GRAIL; and has stock options and is co-founder of Achilles Therapeutics. C.S.
761 holds patents relating to assay technology to detect tumour recurrence
762 (PCT/GB2017/053289); to targeting neoantigens (PCT/EP2016/059401), identifying patent
763 response to immune checkpoint blockade (PCT/EP2016/071471), determining HLA LOH
764 (PCT/GB2018/052004), predicting survival rates of patients with cancer
765 (PCT/GB2020/050221), to treating cancer by targeting Insertion/deletion mutations
766 (PCT/GB2018/051893); identifying insertion/deletion mutation targets
767 (PCT/GB2018/051892); methods for lung cancer detection (PCT/US2017/028013); and
768 identifying responders to cancer treatment (PCT/GB2018/051912).

769

770

771 **Acknowledgements**

772 R.B. is supported by the NIHR BRC at University College London Hospitals. K.L. is funded by
773 the UK Medical Research Council (MR/P014712/1), Rosetrees Trust and Cotswold Trust
774 (A2437) and Cancer Research UK (C69256/A30194). T.B.K.W. was supported by the Francis
775 Crick Institute, which receives its core funding from Cancer Research UK (FC001169), the UK
776 Medical Research Council (FC001169) and the Wellcome Trust (FC001169) as well as the
777 Marie Curie ITN Project PLOIDYNET (FP7-PEOPLE-2013, 607722), Breast Cancer Research
778 Foundation (BCRF), Royal Society Research Professorships Enhancement Award
779 (RP/EA/180007) and the Foulkes Foundation. E.L.L. receives funding from NovoNordisk
780 Foundation (ID 16584). R.R. is supported by Royal Society Research Professorships

781 Enhancement Award (RP/EA/180007). C.M.R is supported by Rosetrees. C.T.H. is supported
782 by the NIHR BRC at University College London Hospitals. M.J.-H. has received funding from
783 Cancer Research UK, National Institute for Health Research, Rosetrees Trust, UKI NETs and
784 NIHR University College London Hospitals Biomedical Research Centre. N.M. is a Sir Henry
785 Dale Fellow, jointly funded by the Wellcome Trust and the Royal Society (Grant Number
786 211179/Z/18/Z) and also receives funding from Cancer Research UK, Rosetrees and the
787 NIHR BRC at University College London Hospitals and the CRUK University College London
788 Experimental Cancer Medicine Centre. C.S. is a Royal Society Napier Research Professor.
789 His work was supported by the Francis Crick Institute, which receives its core funding from
790 Cancer Research UK (FC001169), the UK Medical Research Council (FC001169) and the
791 Wellcome Trust (FC001169). C.S. is funded by Cancer Research UK (TRACERx, PEACE and
792 CRUK Cancer Immunotherapy Catalyst Network), Cancer Research UK Lung Cancer Centre
793 of Excellence, the Rosetrees Trust, Butterfield and Stoneygate Trusts, NovoNordisk
794 Foundation (ID16584), Royal Society Research Professorships Enhancement Award
795 (RP/EA/180007), the NIHR BRC at University College London Hospitals, the CRUK-UCL
796 Centre, Experimental Cancer Medicine Centre and the Breast Cancer Research Foundation
797 (BCRF). This research is supported by a Stand Up To Cancer-LUNGevity-American Lung
798 Association Lung Cancer Interception Dream Team Translational Research Grant (SU2C-
799 AACR-DT23-17). Stand Up To Cancer is a program of the Entertainment Industry Foundation.
800 Research grants are administered by the American Association for Cancer Research, the
801 Scientific Partner of SU2C. C.S. also receives funding from the European Research Council
802 (ERC) under the European Union's Seventh Framework Programme (FP7/2007-2013)
803 Consolidator Grant (FP7-THESEUS-617844), European Commission ITN (FP7-PloidyNet
804 607722), an ERC Advanced Grant (PROTEUS) from the European Research Council under
805 the European Union's Horizon 2020 research and innovation programme (835297) and
806 Chromavision from the European Union's Horizon 2020 research and innovation programme
807 (665233).

808

809 The TRACERx study (Clinicaltrials.gov no: NCT01888601) is sponsored by University College
810 London (UCL/12/0279) and has been approved by an independent Research Ethics
811 Committee (13/LO/1546). TRACERx is funded by Cancer Research UK (C11496/A17786) and
812 coordinated through the Cancer Research UK and UCL Cancer Trials Centre.

813

814 The results published here are based in part on data generated by The Cancer Genome Atlas
815 pilot project established by the NCI and the National Human Genome Research Institute. The
816 data were retrieved through the database of Genotypes and Phenotypes (dbGaP)
817 authorization (accession number phs000178.v9.p8). Information about TCGA and the
818 constituent investigators and institutions of the TCGA research network can be found at
819 <http://cancergenome.nih.gov/>. This project was enabled through access to the MRC eMedLab
820 Medical Bioinformatics infrastructure, supported by the Medical Research Council
821 (MR/L016311/1). In particular, we acknowledge the support of the High-Performance
822 Computing at the Francis Crick Institute as well as the UCL Department of Computer Science
823 Cluster and the support team.

824

825 **The TRACERx Consortium**

826 Charles Swanton (2,4,8), Mariam Jamal-Hanjani (2,8,9), Nicholas McGranahan (1,2), Carlos
827 Martínez-Ruiz (1,2), Robert Bentham (1,2), Kevin Litchfield (2,3), Emilia L Lim (2,4), Crispin T
828 Hiley (2,4), David A Moore (2,7,8), Thomas B K Watkins (4), Rachel Rosenthal (4), Maise Al
829 Bakir (4), Roberto Salgado (5,6), Nicolai J Birkbak (11), Mickael Escudero (11), Aengus
830 Stewart (11), Andrew Rowan (11), Jacki Goldman (11), Peter Van Loo (11), Richard Kevin
831 Stone (11), Tamara Denner (11), Emma Nye (11), Sophia Ward (11), Stefan Boeing (11),
832 Maria Greco (11), Jerome Nicod (11), Clare Puttick (11), Katey Enfield (11), Emma Colliver
833 (11), Brittany Campbell (11), Alexander M Frankell (11), Daniel Cook (11), Mihaela Angelova
834 (11), Alastair Magness (11), Chris Bailey (11), Antonia Toncheva (11), Krijn Dijkstra (11), Judit
835 Kisistok (11), Mateo Sokac (11), Oriol Pich (11), Jonas Demeulemeester (11), Elizabeth
836 Larose Cadieux (11), Carla Castignani (11), Krupa Thakkar (11), Hongchang Fu (11),

837 Takahiro Karasaki (11,12), Othman Al-Sawaf (11,12), Mark S Hill (11,40), Christopher
838 Abbosh (12), Yin Wu (12), Selvaraju Veeriah (12), Robert E Hynds (12), Andrew Georgiou
839 (12), Mariana Werner Sunderland (12), James L Reading (12), Sergio A Quezada (12), Karl
840 S Peggs (12), Teresa Marafioti (12), John A Hartley (12), Helen L Lowe (12), Leah Ensell (12),
841 Victoria Spanswick (12), Angeliki Karamani (12), Dhruva Biswas (12), Stephan Beck (12),
842 Olga Chervova (12), Miljana Tanic (12), Ariana Huebner (12), Michelle Dietzen (12), James
843 RM Black (12), Cristina Naceur-Lombardelli (12), Mita Afroza Akther (12), Haoran Zhai (12),
844 Nnennaya Kanu (12), Simranpreet Summan (12), Francisco Gimeno-Valiente (12), Kezhong
845 Chen (12), Elizabeth Manzano (12), Supreet Kaur Bola (12), Ehsan Ghorani (12), Marc Robert
846 de Massy (12), Elena Hoxha (12), Emine Hatipoglu (12), Benny Chain (12), David R Pearce
847 (12), Javier Herrero (12), Simone Zaccaria (12), Jason Lester (13), Fiona Morgan (14),
848 Malgorzata Kornaszewska (14), Richard Attanoos (14), Haydn Adams (14), Helen Davies
849 (14), Jacqui A Shaw (15), Joan Riley (15), Lindsay Primrose (15), Dean Fennell (15,16),
850 Apostolos Nakas (16), Sridhar Rathinam (16), Rachel Plummer (16), Rebecca Boyles (16),
851 Mohamad Tufail (16), Amrita Bajaj (16), Jan Brozik (16), Keng Ang (16), Mohammed Fiyaz
852 Chowdhry (16), William Monteiro (17), Hilary Marshall (17), Alan Dawson (18), Sara Busacca
853 (18), Domenic Marrone (18), Claire Smith (18), Girija Anand (19), Sajid Khan (19), Gillian Price
854 (20), Mohammed Khalil (20), Keith Kerr (20), Shirley Richardson (20), Heather Cheyne (20),
855 Joy Miller (20), Keith Buchan (20), Mahendran Chetty (20), Sylvie Dubois-Marshall (20), Sara
856 Lock (21), Kayleigh Gilbert (21), Babu Naidu (22), Gerald Langman (22), Hollie Bancroft (22),
857 Salma Kadiri (22), Gary Middleton (22), Madava Djearaman (22), Aya Osman (22), Helen
858 Shackelford (22), Akshay Patel (22), Angela Leek (23), Nicola Totten (23), Jack Davies
859 Hodgkinson (23), Jane Rogan (23), Katrina Moore (23), Rachael Waddington (23), Jane
860 Rogan (23), Yvonne Summers (24), Raffaele Califano (24), Rajesh Shah (24), Piotr Krysiak
861 (24), Kendadai Rammohan (24), Eustace Fontaine (24), Richard Booton (24), Matthew Evison
862 (24), Stuart Moss (24), Juliette Novasio (24), Leena Joseph (24), Paul Bishop (24), Anshuman
863 Chaturvedi (24), Helen Doran (24), Felice Granato (24), Vijay Joshi (24), Elaine Smith (24),
864 Angeles Montero (24), Philip Crosbie (24,25,26), Fiona Blackhall (26,27), Lynsey Priest

865 (26,27), Matthew G Krebs (26,27), Caroline Dive (26,28), Dominic G Rothwell (26,28), Alastair
866 Kerr (26,28), Elaine Kilgour (26,28), Katie Baker (27), Mathew Carter (27), Colin R Lindsay
867 (27), Fabio Gomes (27), Jonathan Tugwood (28), Jackie Pierce (28), Alexandra Clipson (28),
868 Roland Schwarz (29,30), Tom L Kaufmann (31,32), Matthew Huska (33), Zoltan Szallasi (34),
869 Istvan Csabai (35), Miklos Diossy (35), Hugo Aerts (36,37), Charles Fekete (37), Gary Royle
870 (38), Catarina Veiga (38), Marcin Skrzypski (39), David Lawrence (40), Martin Hayward (40),
871 Nikolaos Panagiotopoulos (40), Robert George (40), Davide Patrini (40), Mary Falzon (40),
872 Elaine Borg (40), Reena Khiroya (40), Asia Ahmed (40), Magali Taylor (40), Junaid Choudhary
873 (40), Sam M Janes (40), Martin Forster (40), Tanya Ahmad (40), Siow Ming Lee (40), Neal
874 Navani (40), Dionysis Papadatos-Pastos (40), Marco Scarci (40), Pat Gorman (40), Elisa
875 Bertoja (40), Robert CM Stephens (40), Emilie Martinoni Hoogenboom (40), James W Holding
876 (40), Steve Bandula (40), Ricky Thakrar (40), Radhi Anand (40), Kayalvizhi Selvaraju (40),
877 James Wilson (40), Sonya Hessey (40), Paul Ashford (40), Mansi Shah (40), Marcos Vasquez
878 Duran (40), Mairead MacKenzie (41), Maggie Wilcox (41), Allan Hackshaw (42), Yenting Ngai
879 (42), Abigail Sharp (42), Cristina Rodrigues (42), Oliver Pressey (42), Sean Smith (42), Nicole
880 Gower (42), Harjot Kaur Dhanda (42), Kitty Chan (42), Sonal Chakraborty (42), Christian
881 Ottensmeier (43), Serena Chee (43), Benjamin Johnson (43), Aiman Alzetani (43), Judith
882 Cave (43), Lydia Scarlett (43), Emily Shaw (43), Eric Lim (44), Paulo De Sousa (44), Simon
883 Jordan (44), Alexandra Rice (44), Hilgardt Raubenheimer (44), Harshil Bhayani (44), Morag
884 Hamilton (44), Lyn Ambrose (44), Anand Devaraj (44), Hema Chavan (44), Sofina Begum
885 (44), Silviu I Buderu (44), Daniel Kaniu (44), Mpho Malima (44), Sarah Booth (44), Andrew G
886 Nicholson (44), Nadia Fernandes (44), Christopher Deeley (44), Pratibha Shah (44), Chiara
887 Proli (44), Kelvin Lau (45), Michael Sheaff (45), Peter Schmid (45), Louise Lim (45), John
888 Conibear (45), Madeleine Hewish (46), Sarah Danson (47), Jonathan Bury (47), John
889 Edwards (47), Jennifer Hill (47), Sue Matthews (47), Yota Kitsanta (47), Jagan Rao (47), Sara
890 Tenconi (47), Laura Socci (47), Kim Suvarna (47), Faith Kibutu (47), Patricia Fisher (47), Robin
891 Young (47), Joann Barker (47), Fiona Taylor (47), Kirsty Lloyd (47), Michael Shackcloth (48),
892 Julius Asante-Siaw (48), John Gosney (49), Teresa Light (50), Tracey Horey (50), Peter

893 Russell (50), Dionysis Papadatos-Pastos (50), Kevin G Blyth (51), Craig Dick (51), Andrew
894 Kidd (51), Alan Kirk (52), Mo Asif (52), John Butler (52), Rocco Bilancia (52), Nikos Kostoulas
895 (52), Mathew Thomas (52), Gareth A Wilson (53)
896
897 (11) The Francis Crick Institute, London, United Kingdom
898 (12) University College London Cancer Institute, London, United Kingdom
899 (13) Swansea Bay University Health Board, Swansea, United Kingdom
900 (14) Cardiff & Vale University Health Board, Cardiff, United Kingdom
901 (15) Cancer Research Centre, University of Leicester, Leicester, United Kingdom
902 (16) Leicester University Hospitals, Leicester, United Kingdom
903 (17) National Institute for Health Research Leicester Respiratory Biomedical Research
904 Unit, Leicester, United Kingdom
905 (18) University of Leicester, Leicester, United Kingdom
906 (19) Barnet & Chase Farm Hospitals, Barnet, United Kingdom
907 (20) Aberdeen Royal Infirmary, Aberdeen, United Kingdom
908 (21) The Whittington Hospital NHS Trust, London, United Kingdom
909 (22) University Hospital Birmingham NHS Foundation Trust, Birmingham, United Kingdom
910 (23) Manchester Cancer Research Centre Biobank, Manchester, United Kingdom
911 (24) Wythenshawe Hospital, Manchester University NHS Foundation Trust, Manchester,
912 United Kingdom
913 (25) Division of Infection, Immunity and Respiratory Medicine, University of Manchester,
914 Manchester, UK
915 (26) Cancer Research UK Lung Cancer Centre of Excellence, University of Manchester,
916 Manchester, UK
917 (27) Christie NHS Foundation Trust, Manchester, United Kingdom
918 (28) Cancer Research UK Manchester Institute, University of Manchester, Manchester, UK
919 (29) Berlin Institute for Medical Systems Biology, Max Delbrueck Center for Molecular
920 Medicine, Berlin, Germany

- 921 (30) German Cancer Consortium (DKTK), partner site Berlin, Berlin, Germany
- 922 (31) Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular
923 Medicine in the Helmholtz Association (MDC), Robert-Rössle-Str 10, 13125, Berlin, Germany
- 924 (32) BIFOLD, Berlin Institute for the Foundations of Learning and Data, Berlin, Germany
- 925 (33) Berlin Institute for Medical Systems Biology, Max Delbrueck Center for Molecular
926 Medicine, Berlin, Germany
- 927 (34) Danish Cancer Society Research Center, Copenhagen, Denmark
- 928 (35) Department of Physics of Complex Systems, ELTE Eötvös Loránd University,
929 Budapest, Hungary
- 930 (36) Artificial Intelligence in Medicine (AIM) Program, Mass General Brigham, Harvard
931 Medical School, Boston, Massachusetts, USA
- 932 (37) Radiology and Nuclear Medicine, CARIM & GROW, Maastricht University, Maastricht,
933 The Netherlands
- 934 (38) Department of Medical Physics and Bioengineering, University College London
935 Cancer Institute, London, United Kingdom
- 936 (39) Department of Oncology and Radiotherapy, Medical University of Gdańsk, ul Dębinki
937 7, 80-211, Gdańsk, Poland
- 938 (40) University College London Hospitals, London, United Kingdom
- 939 (41) Independent Cancer Patients Voice, London, United Kingdom
- 940 (42) Cancer Research UK & UCL Cancer Trials Centre, London, United Kingdom
- 941 (43) University Hospital Southampton NHS Foundation Trust, Southampton, United
942 Kingdom
- 943 (44) Royal Brompton and Harefield NHS Foundation Trust, London, United Kingdom
- 944 (45) Barts Health NHS Trust, London, United Kingdom
- 945 (46) Ashford and St Peter's Hospitals NHS Foundation Trust, Chertsey, United Kingdom
- 946 (47) Sheffield Teaching Hospitals NHS Foundation Trust, Sheffield, United Kingdom
- 947 (48) Liverpool Heart and Chest Hospital NHS Foundation Trust, Liverpool, United Kingdom
- 948 (49) Royal Liverpool University Hospital, Liverpool, United Kingdom

949 (50) The Princess Alexandra Hospital NHS Trust, Harlow, United Kingdom

950 (51) NHS Greater Glasgow and Clyde, Glasgow, United Kingdom

951 (52) Golden Jubilee National Hospital, Clydebank, United Kingdom

952 (53) Achilles Therapeutics UK Limited, London, United Kingdom

953

954

955 **References**

956 1. Rosenthal, R. *et al.* Neoantigen-directed immune escape in lung cancer
957 evolution. *Nature* **567**, 479–485 (2019).

958 2. Litchfield, K. *et al.* Meta-analysis of tumor- and T cell-intrinsic mechanisms of
959 sensitization to checkpoint inhibition. *Cell* **184**, 596-614.e14 (2021).

960 3. Robert, C. *et al.* Ipilimumab plus Dacarbazine for Previously Untreated
961 Metastatic Melanoma. *N. Engl. J. Med.* **364**, 2517–2526 (2011).

962 4. Schadendorf, D. *et al.* Pooled analysis of long-term survival data from phase II
963 and phase III trials of ipilimumab in unresectable or metastatic melanoma. *J. Clin.*
964 *Oncol.* **33**, 1889–1894 (2015).

965 5. Goodman, A. M. *et al.* Tumor mutational burden as an independent predictor of
966 response to immunotherapy in diverse cancers. *Mol. Cancer Ther.* **16**, 2598–
967 2608 (2017).

968 6. Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proc. Natl.*
969 *Acad. Sci. U. S. A.* **107**, 16910–16915 (2010).

970 7. Favero, F. *et al.* Sequenza: Allele-specific copy number and mutation profiles
971 from tumor sequencing data. *Ann. Oncol.* **26**, 64–70 (2015).

972 8. Shen, R. & Seshan, V. FACETS: Fraction and Allele-Specific Copy Number
973 Estimates from Tumor Sequencing. *Meml. Sloan-Kettering Cancer Center, Dept.*
974 *Epidemiol. Biostat. Work. Pap. Ser. 1*, 50 (2015).

975 9. Carter, S. L. *et al.* Absolute quantification of somatic DNA alterations in human
976 cancer. *Nat. Biotechnol.* **30**, 413–421 (2012).

- 977 10. López, S. *et al.* Interplay between whole-genome doubling and the accumulation
978 of deleterious alterations in cancer evolution. *Nat. Genet.* **52**, 283–293 (2020).
- 979 11. Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line
980 Encyclopedia. *Nature* **569**, 503–508 (2019).
- 981 12. Levy, E. *et al.* Immune DNA signature of T-cell infiltration in breast tumor
982 exomes. *Sci. Rep.* **6**, 1–10 (2016).
- 983 13. Jamal-Hanjani, M. *et al.* Tracking the Evolution of Non–Small-Cell Lung Cancer.
984 *N. Engl. J. Med.* **376**, 2109–2121 (2017).
- 985 14. Danaher, P. *et al.* Pan-cancer adaptive immune resistance as defined by the
986 Tumor Inflammation Signature (TIS): Results from The Cancer Genome Atlas
987 (TCGA). *J. Immunother. Cancer* **6**, 1–17 (2018).
- 988 15. Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates
989 with markers of immune evasion and with reduced response to immunotherapy.
990 *Science (80-.)*. **355**, (2017).
- 991 16. Aran, D., Hu, Z. & Butte, A. J. xCell: Digitally portraying the tissue cellular
992 heterogeneity landscape. *Genome Biol.* **18**, 1–14 (2017).
- 993 17. Li, T. *et al.* TIMER: A web server for comprehensive analysis of tumor-infiltrating
994 immune cells. *Cancer Res.* **77**, e108–e110 (2017).
- 995 18. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression
996 profiles. *Nat. Methods* **12**, 453–457 (2015).
- 997 19. Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D. E. & Gfeller, D.
998 Simultaneous enumeration of cancer and immune cell types from bulk tumor
999 gene expression data. *Elife* **6**, 1–25 (2017).
- 1000 20. Cancer Genome Atlas Research Network. Comprehensive genomic
1001 characterization of squamous cell lung cancers. *Nature* **489**, 519–525 (2012).
- 1002 21. Cancer Genome Atlas Research Network. Comprehensive molecular profiling of
1003 lung adenocarcinoma. *Nature* **511**, 543–550 (2014).
- 1004 22. Yokoyama, A. *et al.* Age-related remodelling of oesophageal epithelia by

- 1005 mutated cancer drivers. *Nature* **565**, 312–317 (2019).
- 1006 23. Poore, G. D. *et al.* Microbiome analyses of blood and tissues suggest cancer
1007 diagnostic approach. *Nature* **579**, 567–574 (2020).
- 1008 24. Watkins, T. B. K. *et al.* Pervasive chromosomal instability and karyotype order in
1009 tumour evolution. *Nature* **587**, 126–132 (2020).
- 1010 25. Jongsma, M. L. M. *et al.* The SPPL3-Defined Glycosphingolipid Repertoire
1011 Orchestrates HLA Class I-Mediated Immune Responses. *Immunity* **54**, 132-
1012 150.e9 (2021).
- 1013 26. AbdulJabbar, K. *et al.* Geospatial immune variability illuminates differential
1014 evolution of lung adenocarcinoma. *Nat. Med.* **26**, 1054–1062 (2020).
- 1015 27. Schwartz, L. H. *et al.* RECIST 1.1 - Update and clarification: From the RECIST
1016 committee. *Eur. J. Cancer* **62**, 132–137 (2016).
- 1017 28. Shim, J. H. *et al.* HLA-corrected tumor mutation burden and homologous
1018 recombination deficiency for the prediction of response to PD-(L)1 blockade in
1019 advanced non-small-cell lung cancer patients. *Ann. Oncol.* **31**, 902–911 (2020).
- 1020 29. Rizvi, N. A. *et al.* Mutational landscape determines sensitivity to PD-1 blockade
1021 in non-small cell lung cancer. *Science (80-.).* **348**, 124–128 (2015).
- 1022 30. Conforti, F. *et al.* Sex-based dimorphism of anticancer immune response and
1023 molecular mechanisms of immune evasion. *Clin. Cancer Res.*
1024 clincanres.0136.2021 (2021). doi:10.1158/1078-0432.ccr-21-0136
- 1025 31. Capone, I., Marchetti, P., Ascierto, P. A., Malorni, W. & Gabriele, L. Sexual
1026 Dimorphism Of Immune Responses: A new perspective in cancer
1027 immunotherapy. *Front. Immunol.* **9**, 1–8 (2018).
- 1028 32. van der Spek, J., Groenwold, R. H. H., van der Burg, M. & van Montfrans, J. M.
1029 TREC Based Newborn Screening for Severe Combined Immunodeficiency
1030 Disease: A Systematic Review. *J. Clin. Immunol.* **35**, 416–430 (2015).
- 1031
- 1032

1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060

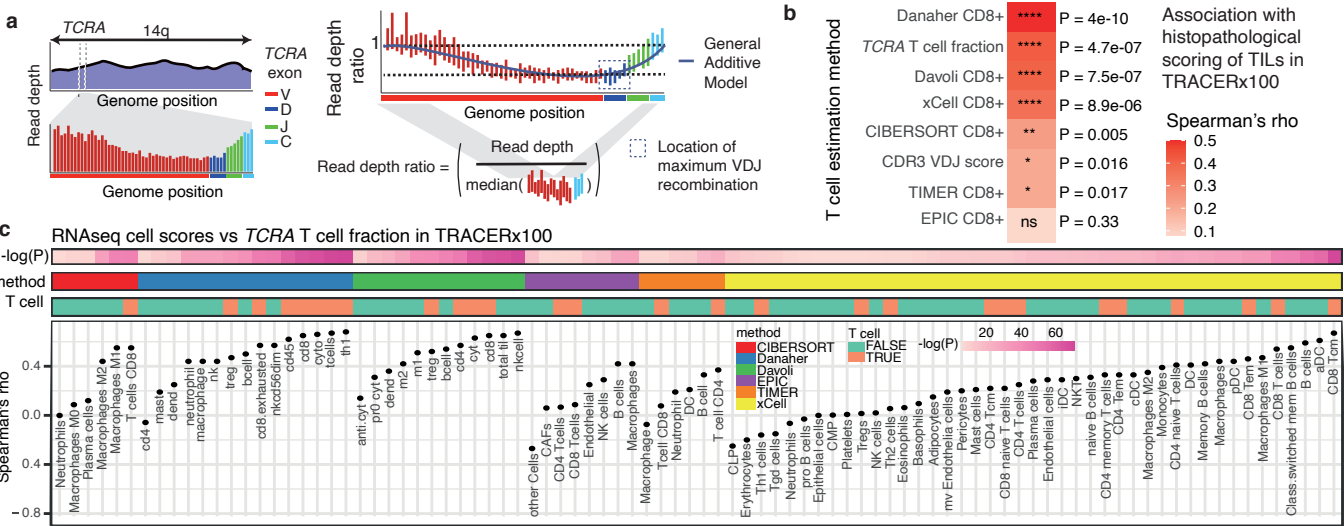
Additional References

33. Kuchenbecker, L. *et al.* IMSEQ-A fast and error aware approach to immunogenetic sequence analysis. *Bioinformatics* **31**, 2963–2971 (2015).
34. Wood, D. E. & Salzberg, S. L. Kraken: Ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15**, (2014).
35. Middleton, G. *et al.* The National Lung Matrix Trial of personalized therapy in lung cancer. *Nature* **583**, 807–812 (2020).
36. Wang, K. *et al.* PennCNV: An integrated hidden Markov model designed for high-resolution copy number variation detection in whole-genome SNP genotyping data. *Genome Res.* **17**, 1665–1674 (2007).
37. Brastianos, P. K. *et al.* Genomic characterization of brain metastases reveals branched evolution and potential therapeutic targets. *Cancer Discov.* **5**, 1164–1177 (2015).
38. Gerlinger, M. *et al.* Genomic architecture and evolution of clear cell renal cell carcinomas defined by multiregion sequencing. *Nat. Genet.* **46**, 225–233 (2014).
39. Gerlinger, M. *et al.* Intratumor Heterogeneity and Branched Evolution Revealed by Multiregion Sequencing. *N. Engl. J. Med.* **366**, 883–892 (2012).
40. Harbst, K. *et al.* Multiregion whole-exome sequencing uncovers the genetic evolution and mutational heterogeneity of early-stage metastatic melanoma. *Cancer Res.* **76**, 4765–4774 (2016).
41. Lamy, P. *et al.* Paired exome analysis reveals clonal evolution and potential therapeutic targets in urothelial carcinoma. *Cancer Res.* **76**, 5894–5906 (2016).
42. Savas, P. *et al.* The Subclonal Architecture of Metastatic Breast Cancer: Results from a Prospective Community-Based Rapid Autopsy Program “CASCADE”. *PLoS Med.* **13**, 1–25 (2016).
43. Suzuki, H. *et al.* Mutational landscape and clonal architecture in grade II and III gliomas. *Nat. Genet.* **47**, 458–468 (2015).
44. Turajlic, S. *et al.* Deterministic Evolutionary Trajectories Influence Primary Tumor

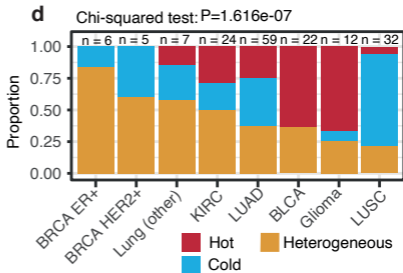
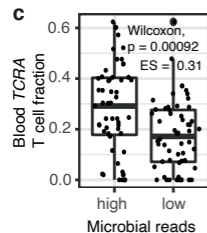
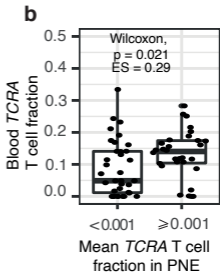
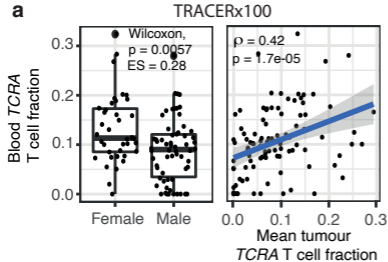
- 1061 Growth: TRACERx Renal. *Cell* **173**, 595-610.e11 (2018).
- 1062 45. Messaoudene, M. *et al.* T-cell bispecific antibodies in node-positive breast
1063 cancer: Novel therapeutic avenue for MHC class I loss variants. *Ann. Oncol.* **30**,
1064 934–944 (2019).
- 1065 46. Snyder, A. *et al.* Genetic Basis for Clinical Response to CTLA-4 Blockade in
1066 Melanoma. *N. Engl. J. Med.* **371**, 2189–2199 (2014).
- 1067 47. Van Allen, E. M. *et al.* Genomic correlates of response to CTLA-4 blockade in
1068 metastatic melanoma. *Science (80-.)*. **352**, 207–212 (2016).
- 1069 48. Hugo, W. *et al.* Genomic and Transcriptomic Features of Response to Anti-PD-
1070 1 Therapy in Metastatic Melanoma. *Cell* **165**, 35–44 (2016).
- 1071 49. Riaz, N. *et al.* Tumor and Microenvironment Evolution during Immunotherapy
1072 with Nivolumab. *Cell* **171**, 934-949.e15 (2017).
- 1073 50. Cristescu, R. *et al.* Pan-tumor genomic biomarkers for PD-1 checkpoint
1074 blockade-based immunotherapy. *Science (80-.)*. **362**, (2018).
- 1075 51. Snyder, A. *et al.* Contribution of systemic and somatic factors to clinical response
1076 and resistance to PD-L1 blockade in urothelial cancer: An exploratory multi-omic
1077 analysis. *PLoS Med.* **14**, 1–24 (2017).
- 1078 52. Mariathasan, S. *et al.* TGF β attenuates tumour response to PD-L1 blockade by
1079 contributing to exclusion of T cells. *Nature* **554**, 544–548 (2018).
- 1080 53. McDermott, D. F. *et al.* Clinical activity and molecular correlates of response to
1081 atezolizumab alone or in combination with bevacizumab versus sunitinib in renal
1082 cell carcinoma. *Nat. Med.* **24**, 749–757 (2018).
- 1083 54. Le, D. T. *et al.* PD-1 blockade in tumors with mismatch repair deficiency. *N. Engl.*
1084 *J. Med.* **372**, 2509–2520 (2015).
- 1085 55. Hendry, S. *et al.* *Assessing Tumor-infiltrating Lymphocytes in Solid Tumors.*
1086 *Advances In Anatomic Pathology* **24**, (2017).
- 1087 56. Denkert, C. *et al.* Standardized evaluation of tumor-infiltrating lymphocytes in
1088 breast cancer: Results of the ring studies of the international immuno-oncology

1089 biomarker working group. *Mod. Pathol.* **29**, 1155–1164 (2016).

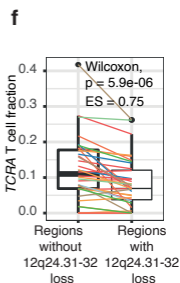
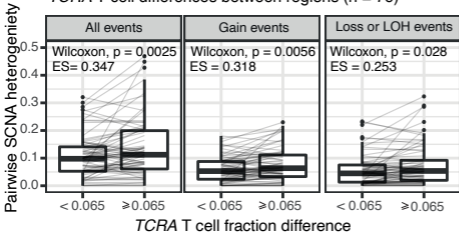
1090



TRACERx100



e Multi-region tumours with both small and large pairwise TCRA T cell differences between regions (n = 76)

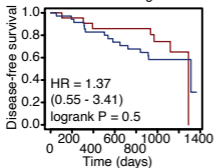


- BLCA
- BRCA ER+
- BRCA HER2+
- ESCA
- Glioma
- KIRC
- LUAD
- Lung (other)
- LUSC
- Lung carcinoma
- SKCM

Decreasing ← Number of immune-cold regions used in threshold → Increasing

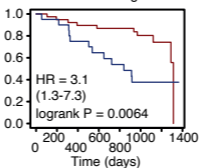
TRACERx100 - LUAD

— < 1 cold regions
— ≥ 1 cold regions



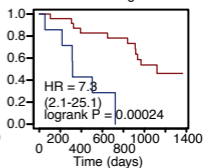
23 20 18 18 17 11 5 0
36 33 29 24 22 16 5 0

— < 2 cold regions
— ≥ 2 cold regions



41 36 34 32 31 22 9 0
20 19 15 12 10 6 1 0

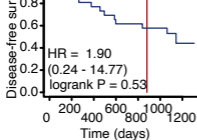
— < 3 cold regions
— ≥ 3 cold regions



25 22 19 18 17 11 3 0
7 6 3 1 0 0 0 0

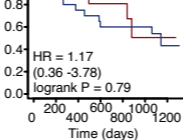
TRACERx100 - LUSC

— < 1 cold regions
— ≥ 1 cold regions



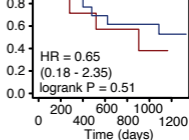
4 4 4 4 2 0 0
28 23 20 16 16 12 5

— < 2 cold regions
— ≥ 2 cold regions



11 9 9 8 6 2 1
21 18 15 12 12 10 4

— < 3 cold regions
— ≥ 3 cold regions



8 6 5 4 3 1 0
14 12 10 8 8 7 4

